

Sequence matching algorithms and pairing of noncoding RNAs

S.K. Nechaev,^{1,2,3} M.V. Tamm,⁴ and O.V. Valba⁵

¹*LPTMS, Université Paris Sud, 91405 Orsay Cedex, France*

²*P.N. Lebedev Physical Institute of the Russian Academy of Sciences, 119991, Moscow, Russia*

³*J.-V. Poncelet Laboratory, Independent University, 119002, Moscow, Russia*

⁴*Physics Department, Moscow State University, 119992, Moscow, Russia*

⁵*Moscow Institute of Physics and Technology, 141700, Dolgoprudny, Russia*

(Dated: November 12, 2010)

A new statistical method of alignment of two heteropolymers which can form hierarchical cloverleaf-like secondary structures is proposed. This offers a new constructive algorithm for quantitative determination of binding free energy of two noncoding RNAs with arbitrary primary sequences. The alignment of ncRNAs differs from the complete alignment of two RNA sequences: in ncRNA case we align only the sequences of nucleotides which constitute pairs between two different RNAs, while the secondary structure of each RNA comes into play only by the combinatorial factors affecting the entropic contribution of each molecule to the total cost function. The proposed algorithm is based on two observations: i) the standard alignment problem is considered as a zero-temperature limit of a more general statistical problem of binding of two associating heteropolymer chains; ii) this last problem is generalized onto the sequences with hierarchical cloverleaf-like structures (i.e. of RNA-type). Taking zero-temperature limit at the very end we arrive at the desired “cost function” of the system with account for entropy of side cactus-like loops. Moreover, we have demonstrated in detail how our algorithm enables to solve the “structure recovery” problem. Namely, we can predict in zero-temperature limit the cloverleaf-like (i.e. secondary) structure of interacting ncRNAs by knowing only their primary sequences.

Contents

I. Introduction	2
A. Binding of noncoding RNAs	2
B. Noncoding RNAs as particular class of associating heteropolymers	3
C. Pairing vs alignment	4

II. Theoretical background	5
A. Alignment of linear sequences	5
B. Matching vs pairing of two random linear heteropolymers	6
C. Matching vs pairing of two random RNA-type heteropolymers	8
III. Matching algorithm for two noncoding RNAs	11
A. Matching of linear sequences	11
B. Matching of RNA-type sequences	12
IV. Structure recovery	13
A. Finding the Longest Common Subsequence for linear chains	13
B. Finding the secondary structure for interacting RNA-like chains	14
V. Conclusion	18
A. Temperature dependence of the free energy	20
B. Statistical analysis of a pair of random sequences	21
References	22

I. INTRODUCTION

A. Binding of noncoding RNAs

According to a common definition, the noncoding RNA (ncRNA) is an RNA molecule that is not translated into a protein. The ncRNAs either regulate the gene expression directly, for example by occupying the ribosome binding site, or indirectly providing RNA targeting specificity for a protein-based regulatory mechanism [1]. The class of ncRNAs spreads on regulatory and functional RNAs. In modern classification the term “noncoding RNA” is basically attributed to eucariotic RNAs, sometimes called also “small nonmessenger RNAs”. In general, regulatory RNAs act in the cell by one of two basic mechanisms: by base-pairing interactions with other nucleic acids, or by binding to proteins [2]. The base pairing with target molecules constitutes the typical mechanism, by which the ncRNA regulates the gene expression. The base pairing is subdivided in two classes depending on their locations: *cis*-encoded ncRNAs are placed at the same genetic location but on the strand opposite to

the target RNA, and *trans*-encoded ncRNAs are placed at a chromosomal location distinct from the target RNA.

It should be noted however that the direct gene regulation of ncRNA via specific base pairing is not completely understood. One of few known examples concerns the participation of the *Xist* gene in *X*-inactivation [3]. *X*-inactivation (also called *lyonization*) is a process by which one of two copies of *X* chromosome present in female mammals is inactivated. The inactive *X* chromosome is silenced by packaging into transcriptionally inactive heterochromatin. The *Xist* gene exhibits properties of the *X*-inactivation center and *Xist* ncRNA becomes localized close to the autosome into which the gene is integrated [3].

Since base-pairing of noncoding and target RNAs plays such important biological role, it is worth estimating theoretically the binding free energy of the ncRNA–target RNA complex by knowing the primary structures of each macromolecule. This problem resembles the alignment of two RNA sequences with one principal difference: in ncRNA case we align only the sequences of nucleotides which constitute pairs between two RNAs, while the secondary structure of each RNA comes into play only by the combinatorial factors affecting the entropic contribution to the total cost function.

One of the key problems in computational ncRNA genefinding is to predict RNA transcript initiation, termination, and processing. However, accurate prediction of even simple transcription units remains an open question – see, for example, the minireview [4].

In brief, the main goal of this work consists in developing a constructive method to build a “cost function”, which characterizes matching (alignment) of two noncoding RNAs with arbitrary primary sequences.

B. Noncoding RNAs as particular class of associating heteropolymers

To put problem of alignment of ncRNAs into the context of statistical mechanics, it seems desirable to extract the basic features of ncRNAs which would play the major role in our analysis. The ncRNAs are the particular examples of a wide class of so-called “associating” heteropolymers.

Generally, associating polymers, besides the strong covalent interactions responsible for the frozen primary sequence of monomer units, are capable of forming additional weaker reversible temperature-dependent (i.e. “thermoreversible”) bonds between different monomers. Many biologically important macromolecules, like proteins and nucleic acids, belong to the class of associating polymers [5].

For associating polymers the variety of possible thermodynamic states and ternary structures is determined by the interplay between the following three major factors: i) the energy

gain due to the direct “pairing”, i.e. formation of thermoreversible contacts; ii) the combinatoric entropy due to the choice of which particular monomers (among those able to participate in bonds formation) do actually create bonds; iii) the loss of conformational entropy of the polymer chain due to pairing (and in particular, the entropic penalty of loop creation between two paired monomers).

Among a variety of macromolecular systems with thermoreversible pairing we pay a special attention to a class of RNA-like polymers. These polymers are distinguished from other biologically active associating polymers, such as, for instance, proteins, by a capability of forming hierarchical “cloverleaf-like” (or “cactus-like”) secondary structures. The formation of a thermoreversible contact between two distant bonds in a RNA molecule (or in a single-stranded DNA) imposes a nonlocal constraint on a number of unpaired possible conformations: all bonds in a RNA chain are known to be arranged in a way to allow only hierarchical cactus-like folded conformations topologically isomorphic to a tree. The pairs of bonds, which do not obey such a structure are called “pseudoknots”. In most cases they are forbidden for RNA molecules. We shall accept the absence of pseudoknots as a matter of fact. Let us note however that in the work [6] the dynamic programming algorithm has been developed for predicting optimal RNA secondary structure, including pseudoknots.

Being formulated in statistical terms, the main goal of our work can be rephrased as follows. We propose a new efficient and statistically justified algorithm for the determination of the binding free energy of any two primary heteropolymer sequences under the supposition that each sequence can form a hierarchical cactus-like secondary structure, typical for RNA molecules.

C. Pairing vs alignment

Let us reveal the similarities and differences between computations of the free energy of associating heteropolymer complexes and standard matching algorithms.

The matching (or “alignment”) problem, even for linear structures is one of the key tasks of computational evolutionary biology. In particular, one of the most important applications of Longest Common Subsequence (LCS) search in linear structures is a quantitative definition of a “closeness” of two DNA sequences. Such a comparison provides information about how far, in evolutionary terms, two genes of one parent have deviated from each other. Also, when a new DNA molecule is sequenced *in vitro*, it is important to know whether it is really new or it is similar to already existing molecules. This is achieved quantitatively by measuring the LCS of the new molecule with other ones available from databases.

The task of the present work consists of extending the statistical approach developed for alignment of linear sequences to the computation of pairing free energy of two RNA-type

structures. The target object of our approach is a ground state free energy as complexes ncRNA – target RNA, or ncRNA – DNA.

II. THEORETICAL BACKGROUND

A. Alignment of linear sequences

Recall that the problem of finding the LCS of a pair of linear sequences drawn from the alphabet of c letters is formulated as follows. Consider two sequences $\alpha = \{\alpha_1, \alpha_2, \dots, \alpha_m\}$ (of length m) and $\beta = \{\beta_1, \beta_2, \dots, \beta_n\}$ (of length n). For example, let α and β be two random sequences of $c = 4$ base pairs A, C, G, T of a DNA molecule, e.g., $\alpha = \{A, C, G, C, T, A, C\}$ with $m = 6$ and $\beta = \{C, T, G, A, C\}$ with $n = 5$. Any subsequence of α (or β) is an ordered sublist of α (and of β) entries which need not to be consecutive, e.g, it could be $\{C, G, T, C\}$, but not $\{T, G, C\}$. A common subsequence of two sequences α and β is a subsequence of both of them. For example, the subsequence $\{C, G, A, C\}$ is a common subsequence of both α and β . There are many possible common subsequences of a pair of initial sequences. The aim of the LCS problem is to find the longest of them. This problem and its variants have been widely studied in biology [7–10], computer science [11–14], probability theory [16–21] and more recently in statistical physics [15, 22–24].

The basis of dynamic programming algorithms for comparing genetic sequences has been formulated for the first time in [25] (see also [26]). In general setting this algorithm takes into account the number of perfect matches and the difference between mismatches and gaps. Being formulated in statistical terms, it consists in constructing the “cost function”, F , having a meaning of an energy (see, for example [27, 28] for details)

$$F = N_{\text{match}} + \mu N_{\text{mis}} + \delta N_{\text{gap}} \quad (1)$$

In Eq.(1) N_{match} , N_{mis} and N_{gap} are correspondingly the numbers of matches, mismatches and gaps in a given pair of sequences, and μ and δ are respectively the energies of mismatches and gaps. Without the loss of generality, the energy of matches can be always set to 1. Besides Eq.(1) we have an obvious conservation law

$$n + m = 2N_{\text{match}} + 2N_{\text{mis}} + N_{\text{gap}} \quad (2)$$

which allows one to exclude N_{gap} from Eq.(1) and rewrite it as follows:

$$F = N_{\text{match}} + \mu N_{\text{mis}} + \delta(n + m - 2N_{\text{match}} - 2N_{\text{mis}}) = (1 - 2\delta)N_{\text{match}} + (\mu - 2\delta)N_{\text{mis}} + \text{const} \quad (3)$$

In Eq.(3) the irrelevant constant $\delta(n + m)$ can be dropped out.

Adopting $(1 - 2\delta)$ as a unit of energy, we arrive at the following expression

$$\tilde{F} = N_{\text{match}} + \gamma N_{\text{mis}} \quad (4)$$

where

$$\gamma = \frac{\mu - 2\delta}{1 - 2\delta}, \quad (5)$$

and $\gamma \leq 1$ by definition. The interesting region is $0 \leq \gamma \leq 1$, otherwise there are no mismatches at all in the ground state (i.e., there is no difference between $\gamma = 0$, which corresponds to simplest version of the LCS problem, and $\gamma < 0$).

It is known [27, 28] that the maximal cost function

$$\tilde{F}^{\max} = \max [N_{\text{match}} + \gamma N_{\text{mis}}] \quad (6)$$

can be computed recursively using the “dynamic programming”

$$\tilde{F}_{m,n}^{\max} = \max \left[\tilde{F}_{m-1,n}^{\max}, \tilde{F}_{m,n-1}^{\max}, \tilde{F}_{m-1,n-1}^{\max} + \zeta_{m,n} \right] \quad (7)$$

with

$$\zeta_{m,n} = \begin{cases} 1 & \text{in case of match} \\ \gamma & \text{in case of mismatch} \end{cases} \quad (8)$$

In our previous studies of matching statistics in *linear* sequences we have shown in [29] that properly normalized asymptotic distribution of the LCS in a somewhat simplified version of the problem, known in literature as a “Bernoulli model”, is given by the so-called Tracy–Widom distribution, first derived for the distribution of the highest eigenvalues of random matrices belonging to the Gaussian ensemble [30, 31].

B. Matching vs pairing of two random linear heteropolymers

Consider the auxiliary statistical model describing the formation of a complex of two heteropolymer linear chains with arbitrary primary sequences. Let these chains be of lengths $L_1 = m\ell$ and $L_2 = n\ell$ correspondingly. In what follows we shall measure the lengths of the chains in number of monomers, m and n , supposing that the size of an elementary unit, ℓ , is equal to 1. Every monomer can be chosen from a set of c different types A, B, C, D, Monomers of the first chain could form saturating reversible bonds with monomers of the second chain. The term “saturating” means that any monomer can form a bond with at most one monomer of the other chain. The bonds between similar types (like A–A, B–B, C–C, etc.) have the attraction energy u and are called below “matches”, while the bonds between different types (like A–B, A–D, B–D, etc.) have the attraction energy v and are called “mismatches” [38]. Suppose also that some parts of the chains can form loops. These loops obviously produce “gaps” since the monomers inside the loops of one chain have no matching (or mismatching) counterparts in the other chain. Schematically a particular configuration of the system under consideration for $c = 2$ is shown in Fig.1.

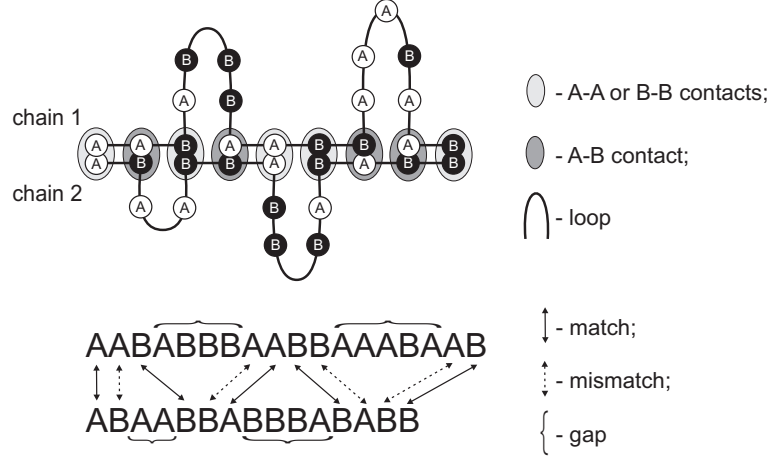


Figure 1: Schematic picture of a complex of two random linear heteropolymer chains with two types of letters ($c = 2$).

Our aim is to compute the free energy of the described model at sufficiently low temperatures under the supposition that the entropic contribution of the loop formation is negligible compared to the energetic part of the direct interactions between chain monomers. Let $G_{m,n}$ be the partition function of such a complex; $G_{m,n}$ is the sum over all possible arrangements of bonds. In the low-temperature limit we can write $G_{m,n}$ recursively:

$$\begin{cases} G_{m,n} = 1 + \sum_{i,j=1}^{m,n} \beta_{i,j} G_{i-1,j-1} \\ G_{m,0} = 1; G_{0,n} = 1; G_{0,0} = 1 \end{cases} \quad (9)$$

The meaning of the equation (9) is straightforward. Starting from, say, the left ends of the chains shown in Fig.1 we find the first actually existing contact between the monomers i (of the first chain) and j (of the second chain) and sum over all possible arrangements of this first contact. The first term "1" in (9) means that we have not found any contact at all. The entries $\beta_{i,j}$ ($1 \leq i \leq m$, $1 \leq j \leq n$) are the statistical weights of the bonds which are encoded in a contact map $\{\beta\}$:

$$\beta_{m,n} = \begin{cases} \beta^+ \equiv e^{u/T} & \text{monomers } i \text{ and } j \text{ match} \\ \beta^- \equiv e^{v/T} & \text{monomers } i \text{ and } j \text{ do not match} \end{cases} \quad (10)$$

The straightforward computation shows that the partition function $G_{m,n}$ (9) obeys the following exact local recursion

$$G_{m,n} = G_{m-1,n} + G_{m,n-1} + (\beta_{m,n} - 1) G_{m-1,n-1} \quad (11)$$

Note that if $\beta_{i,j} = 2$ for all $1 \leq i \leq m$ and $1 \leq j \leq n$, the recursion relation (11) generates the so-called Delannoy numbers [33].

Let us point out that since we are working at finite temperatures, the account for “loop factors” is desirable. Under the “loop factor” we understand the entropic contribution to the free energy of the entire system coming from the fluctuations of parts of heteropolymer chains between successive contacts. Obviously, in the zero-temperature limit these fluctuations vanish.

Write the partition function $G_{m,n}$ as $G_{m,n} = \exp\{F_{m,n}/T\}$, where $-F_{m,n}$ and T are the free energy and the temperature of the complex of two heterogeneous chains of lengths m and n . Considering the $T \rightarrow 0$ limit of the equation (11), we get

$$F_{m,n} = \lim_{T \rightarrow 0} T \ln \left(e^{F_{m-1,n}/T} + e^{F_{m,n-1}/T} + (\beta_{m,n} - 1) e^{F_{m-1,n-1}/T} \right) \quad (12)$$

which can be regarded as an equation for the ground state energy of a chain. The expression (12) reads

$$F_{m,n} = \max [F_{m-1,n}, F_{m,n-1}, F_{m-1,n-1} + \eta_{m,n}] \quad (13)$$

where

$$\eta_{m,n} = T \ln(\beta_{m,n} - 1) = \begin{cases} \eta^+ = T \ln(e^{u/T} - 1) & \text{match} \\ \eta^- = T \ln(e^{v/T} - 1) & \text{mismatch} \end{cases} \quad (14)$$

Indeed, the ground state energy (13) may correspond either: (i) to the last two monomers connected, then the ground state energy equals $\tilde{F}_{m-1,n-1}^{\max} + \zeta_{M,N}$, or (ii) to the unconnected end monomer of the first (or second) chain, then the ground state energy is $\tilde{F}_{m,n-1}^{\max}$ (or $\tilde{F}_{m-1,n}^{\max}$).

Taking η^+ as the unit of the energy, rewrite (13) in a form identical to the dynamic programming equation (7):

$$\tilde{F}_{m,n} = \max [\tilde{F}_{m-1,n}, \tilde{F}_{m,n-1}, \tilde{F}_{m-1,n-1} + \tilde{\eta}_{m,n}] \quad (15)$$

with

$$\tilde{\eta}_{m,n} = \begin{cases} 1 & \text{in case of match} \\ a = \frac{\eta^-}{\eta^+} & \text{in case of mismatch} \end{cases} \quad (16)$$

(compare to (8)). In the low-temperature limit the parameter a has simple expression in terms of coupling constants u and v :

$$a = \frac{\eta^-}{\eta^+} = \frac{\ln(e^{v/T} - 1)}{\ln(e^{u/T} - 1)} \Big|_{T \rightarrow 0} = \frac{v}{u} \quad (17)$$

The initial conditions for $\tilde{F}_{m,n}$ are transformed into $\tilde{F}_{0,n} = \tilde{F}_{n,0} = \tilde{F}_{0,0} = 0$.

C. Matching vs pairing of two random RNA-type heteropolymers

Having the applications to RNA molecules in mind, assume that the structures formed by thermoreversible bonds of each chain are always of a cactus-like type, as shown in Fig.2a.

It means that we restrict ourselves to the situation in which the chain conformations with "pseudoknots" shown in Fig.2b are prohibited. The difference between allowed and not allowed structures becomes more transparent, being redrawn in the following way. Represent a polymer under consideration as a straight line with active monomers situated along it in the natural order, and depict the bonds by dashed arcs connecting the corresponding monomers. Now, the absence of pseudoknots means the absence of intersection of the arcs – see the Fig.2c,d.

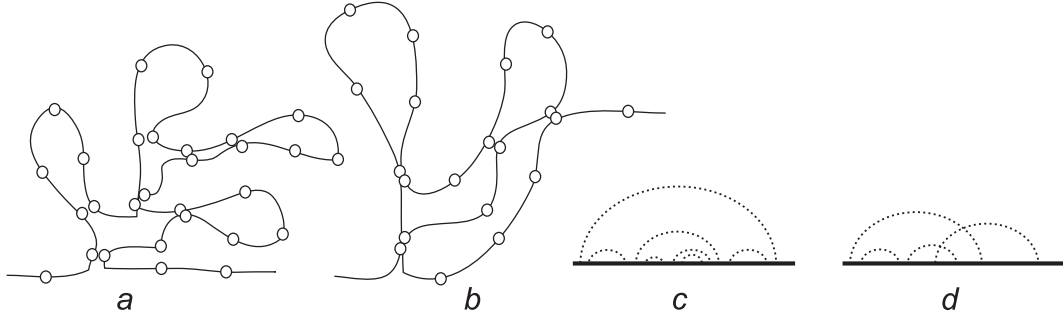


Figure 2: (a,b): Schematic picture of allowed (a) cactus-like and prohibited (b) pseudoknot configurations of the bonds; (c,d): Arc diagrams corresponding respectively to configurations (a) and (b) (note the intersection of arcs in (d)).

We assume for simplicity, that except pseudoknots, all other bond configurations are allowed. This means, in particular, that at the moment we do not require any minimal loop length, as well as we do not yet take into account the cooperativity effect [39]. These assumptions are known to be false for real RNA molecules (for example, there are no loops shorter than 3 monomers in RNA chains [34]). However, one can speculate that (see, for example, [35]) if the links of the chain are considered as renormalized quasi-monomers consisting of several "bare" units, these assumptions seem to be plausible. Nevertheless, in the last Section we study in detail the effect of minimal loop length on the structure formation.

Let us remind that one of the main goals in this work consists in developing an algorithm for the computation of the cost function, which characterizes the similarity of two RNA-type random sequences. To succeed, we should incorporate in the conventional cost function discussed above the contribution coming from the entropy of different rearrangements of cactus-like conformations typical for RNA's. It is not obvious how to do that directly in the frameworks of the dynamic programming approach formalized in the recursion relation (7). To proceed, we exploit the idea (formulated for the first time in [32]), which consists of two consecutive steps:

1. First of all, we reformulate the recursion relation (7) in terms natural for statistical mechanical consideration and show that (7) can be regarded as a relation for the free

energy of some statistical model describing the formation of a complex of two random heteropolymer linear chains in a zero-temperature limit;

2. Secondly, we take into account the possibility for random heteropolymer chains to form complex spatial cactus-like structures and write the corresponding recursion relations for the *partition function* (but not for the free energy) at some temperature T not obliged to be zero. By taking the limit $T \rightarrow 0$ at the very end we arrive at the desired cost function.

The generic partition function $G_{m,n}$ of a complex of two heteropolymers, where each of chains can form a cactus-like structure, shown in the Fig.2a, can be written in the form similar to (9):

$$\begin{cases} G_{m,n} = g_m^{(1)} g_n^{(2)} + \sum_{i,j=1}^{m,n} \beta_{i,j} G_{i-1,j-1} g_{m-i}^{(1)} g_{n-j}^{(2)} \\ G_{m,0} = g_m^{(1)}; \quad G_{0,n} = g_n^{(2)}; \quad G_{0,0} = 1 \end{cases} \quad (18)$$

where $g_n^{(1)}$ and $g_m^{(2)}$ are the partition functions of individual chains. They satisfy the selfconsistent Dyson-type equation [34, 36, 37]

$$g_n^{(1)} = 1 + \sum_{i=1}^{n-1} \sum_{j=i+1+\ell}^n \beta'_{i,j} \frac{g_{j-i-1}^{(1)}}{(j-i-1)^\alpha} g_{n-j}^{(1)}; \quad g_0^{(1)} = 1 \quad (19)$$

(the same equation should be written for $g_m^{(2)}$). The Boltzmann weights $\beta'_{i,j}$ are the constants of self-association, which are, similarly to $\beta_{m,n}$, variables encoded by some contact map and the denominator describes the contribution of the entropic “loop factor”. The value $\alpha = 3/2$ (considered throughout our paper) corresponds to the loop factor of ideal chains. The summation over j running from $i+1+\ell$ till n ensures the absence of loops of lengths smaller than $\ell = 3$ monomers. The equation (19) is schematically depicted in the Fig.3.



Figure 3: Diagrammatic form of the Dyson-type equation for the partition function of an individual chain g_n having cactus-like topology.

Equations (18)–(19) constitute the analytical basis of our numerical studies and these equations are considered as a replacement of the dynamic programming algorithm for matching of sequences with RNA-type architecture.

III. MATCHING ALGORITHM FOR TWO NONCODING RNAs

In this Section we describe an algorithm for computing the binding free energy (which plays a role of the cost function) for the pair of two noncoding RNAs. Let us remind that in ncRNA case we align only the sequences of nucleotides which constitute pairs between two different RNAs and the cactus-like secondary structure of each RNA contributes to the total cost function by corresponding entropic factors.

Extrapolating the free energy of linear sequences to zero temperature we recover (for linear sequences only) the well-known standard dynamic programming algorithm described in (15)–(17). For cactus-like structures our algorithm is not reduced (even at zero temperature) to any local recursive scheme.

The readers who are not interested in the details of the mathematical background discussed at length of the Section II, can regard the results of the current Section as a self-contained prescription for the computation of the desired cost function.

For clarity we formulate the sequential steps of our algorithm keeping in mind two trial sequences of nucleotides of lengths m and n with $m = n = 75$. These sequences are depicted in the Fig.4. These sequences will be aligned in two ways being considered as linear and cactus-like (“RNA-like”). The free energy (i.e. the cost function) of two sequences of total lengths m and n is

$$F_{m,n} = T \ln G_{m,n} \quad (20)$$

AUCGAUGUAGGGUACACGGGCUUAUGUUACGACGAGAUGUCUUGUUCGAUCAUGCGCUUCCGCGGAGAGUGGAAA - s1
AGUUGCACCGCCAGACUACUUAACUAAACGUCGGCCAAGACAAUUCGCAUCGACCUAGUUAGCACGCACCAUCGA - s2

Figure 4: Two trial sequences of $m = n = 75$ nucleotides.

A. Matching of linear sequences

Suppose for the time being that both sequences in Fig.4 are linear. Construct the matrix G whose elements $G_{i,j}$ ($1 \leq i \leq m; 1 \leq j \leq n$) are the partition functions satisfying the relation (10)–(11) with the boundary conditions $G_{m,0} = G_{0,n} = G_{0,0} = 1$ (see (9)). The matrix element $G_{i,j}$ defines matching of i first nucleotides of the 1st sequence with j first nucleotides of the 2nd one. The effective energy of two complimentary nucleotides in the (10) is $u = 1$, while for non-complimentary ones is $v = 0$. It is easy to see from (11) that the search of $G_{m,n}$ can be completed in polynomial time $\sim O(mn)$. At $T \rightarrow 0$ we recover the standard dynamic programming algorithm [25, 26] (see (7)).

B. Matching of RNA-type sequences

Suppose now that both sequences in Fig.4 can form hierarchical cactus-like (i.e. “RNA-type”) structures. The computation of the free energy of the complex built by the pair of RNA-type sequences can be accomplished in two sequential steps:

- Compute the matrices $g^{(1)}$ and $g^{(2)}$ (of sizes $m \times m$ and $n \times n$) of statistical weights of 1st and 2nd sequences separately. Rewrite (19) as

$$g_{i,j}^{(a)} = 1 + \sum_{r=i}^{j-1} \sum_{s=i+1+\ell}^j \beta'_{r,s} \frac{g_{r+1,s-1}^{(a)}}{(s-r-1)^\alpha} g_{s+1,j}^{(a)}; \quad g_{i,i}^{(a)} = 1 \quad (21)$$

where $g_{i,j}^{(a)}$ is the statistical weight of the loop from the nucleotide i till the nucleotide j in the 1st ($a = 1$) or 2nd ($a = 2$) sequence. For each $a = 1, 2$ the systems of equations (21) are quadratic in $g_{i,j}^{(a)}$ and can be solve recursively. The boundary conditions together with the recursion scheme (21) uniquely define the elements $g_{i,i+1}^{(a)}$. Knowing $g_{i,i+1}^{(a)}$ and applying (21) again, we compute $g_{i,i+2}^{(a)}$. The elements $g_{i,j}^{(a)}$ with $i > j$ are set equal to zero. The free energy (the cost function) of the hierarchical cactus-like structure is defined by (20).

- Knowing the matrices $g^{(1)}$ and $g^{(2)}$ find the elements $G_{i,j}$ of the matrix G by solving (18).

The ground state free energy $F_0 \equiv F(T = 0)$ (i.e. the binding free energy at zero’s temperature) for cactus-like structures can be explicitly computed by extending the approach, developed in Section II C. The zero-temperature free energies $F_{m,n}$ of branching structures read (compare to Eqs.(15)–(16)):

$$F_{m,n} = \max_{\substack{i=1,\dots,m \\ j=1,\dots,n}} \left[f_{1,m}^{(1)} + f_{1,n}^{(2)}, Q_{i,j} \right] \quad (22)$$

where $f_{i,j}^{(a)} = T \ln g_{i,j}^{(a)}$ ($a = 1, 2$) are the free energies of individual subsequences from the nucleotide i till the nucleotide j , and $Q_{i,j}$ is the zero-temperature limit of the (i, j) term in Eq.(18):

$$Q_{i,j} = F_{i-1,j-1} + f_{i+1,m}^{(1)} + f_{j+1,n}^{(2)} + \tilde{\eta}_{i,j} \quad (23)$$

At $T = 0$ one can write

$$f_{i,j}^{(a)} = \max_{\substack{r=1,\dots,i \\ s=i+1+\ell,\dots,j}} \left[f_{r+1,s-1}^{(a)} + f_{s+1,j}^{(a)} + \tilde{\eta}'_{r,s} \right] \quad (24)$$

The values $\tilde{\eta}_{i,j}$ define the matching constants of linear sequences (as in (16)), while $\tilde{\eta}_{i,j}^{(a)}$ are the matching constants in each separate sequence.

The boundary conditions for the ground state free energy follow from the boundary conditions of the partition function (18):

$$\begin{cases} F_{0,0} = 0; \\ F_{i,0} = f_{1,i}^{(1)}; \quad 1 \leq i \leq m \\ F_{0,j} = f_{1,j}^{(2)}; \quad 1 \leq j \leq n \end{cases} \quad (25)$$

Thus, to compute the ground state free energy of the complex of two RNA-like sequences, we should first reconstruct the matrices $f^{(1)}$ and $f^{(2)}$ of individual chains by applying Eq.(24) and then find the matrix F by using Eq.(22). The boundary conditions (25) together with Eq.(23) allow us to compute the elements of the matrix Q for $m = 1$ and any n . Knowing the corresponding matrix Q we define the elements $F_{1,j}$ ($1 \leq j \leq n$) of the free energy matrix by using Eq.(22). Then we proceed recursively and determine the matrix Q for $m = 2$ and any n , compute again $F_{2,j}$ ($1 \leq j \leq n$) etc.

For the sequence depicted in Fig.4 we have found the following values of the ground state free energies:

$$\begin{cases} F_1(T = 0) = 48 \quad \text{for linear structure} \\ F_c(T = 0) = 51 \quad \text{for cactus-like structure with } \alpha = 3/2 \text{ and } \ell > 3 \end{cases}$$

The obtained values coincide with the total number of complimentary pairs in formed (linear or cactus-like) structures. The discussion of the temperature behavior of the free energy, $F(T)$, is given in the Appendix A.

IV. STRUCTURE RECOVERY

In this Section we describe the implementation of the structure recovery algorithm for linear and cactus-like structures by the corresponding matrices of free energies F at zero temperature. Let us point out that due to the degeneration mentioned above, the restored sequence is one among the ensemble of sequences with the same free energy.

A. Finding the Longest Common Subsequence for linear chains

Sequence matching problem for linear structures consists in finding the longest common subsequence (possible with gaps) of two given sequences of nucleotides. Let us demonstrate on simple example how the algorithm works. Consider two sequences of $m = n = 6$ nucleotides and construct the incidence matrix η with $\eta_{i,j} = 1$ if monomers i of the 1st sequence and j of the second one match each other, and $\eta_{i,j} = 0$ otherwise – see Fig.5a. In Fig.5b

		G	C	G	G	A	A
	C	1	0	1	1	0	0
	G	0	0	1	0	0	0
$\eta =$	U	0	0	0	0	1	1
	U	0	0	0	0	1	1
	C	1	0	1	1	0	0
	C	1	0	1	1	0	0

(a)

		G	C	G	G	A	A
	C	1	1	1	1	1	1
	G	1	2	2	2	2	2
$F =$	U	1	2	2	2	3	3
	U	1	2	2	2	3	4
	C	1	2	3	3	3	4
	C	1	2	3	4	4	4

(b)

Figure 5: (a) Incidence matrix η , (b) ground state free energy matrix F .

we have shown the matrix of ground state free energies, F , computed via the recursion algorithm (15)–(16).

In order to see which nucleotides form links, let us proceed as follows. Take the element $F_{i,j}$ of the matrix F and compare its value to the values of three neighboring matrix elements $F_{i-1,j-1}$, $F_{i-1,j}$, $F_{i,j-1}$. Now we take the following decisions:

- If $F_{i-1,j-1} = \max[F_{i-1,j-1}, F_{i-1,j}, F_{i,j-1}]$ then i of the 1st sequence is linked to j of the 2nd one;
- If $F_{i-1,j} = \max[F_{i-1,j-1}, F_{i-1,j}, F_{i,j-1}]$ then we skip the element i in the 1st sequence;
- If $F_{i,j-1} = \max[F_{i-1,j-1}, F_{i-1,j}, F_{i,j-1}]$ then we skip the element j in the 2nd sequence.

This procedure begins with the element $F_{m,n}$.

This prescription for computing the matrix of ground state free energies shown in Fig.5b gives (due to degeneration) many sequences with the same value of the free energy. Two possible realizations are depicted in Fig.6.

B. Finding the secondary structure for interacting RNA-like chains

The structure recovery for the chains with cactus-like structures is much more involved problem, however it can also be described recursively. In this case the algorithm consists of the following successive steps:

- Begin with the element $F_{m,n}$ and use (22). If $F_{m,n} > f_{1,m}^{(1)} + f_{1,n}^{(2)}$ we consider the matrix Q (Eq.(23)) and chose the maximal element $Q_{p,q}$ of the matrix Q which corresponds to pairing between the nucleotide p of the 1st sequence and nucleotide q of the second one;

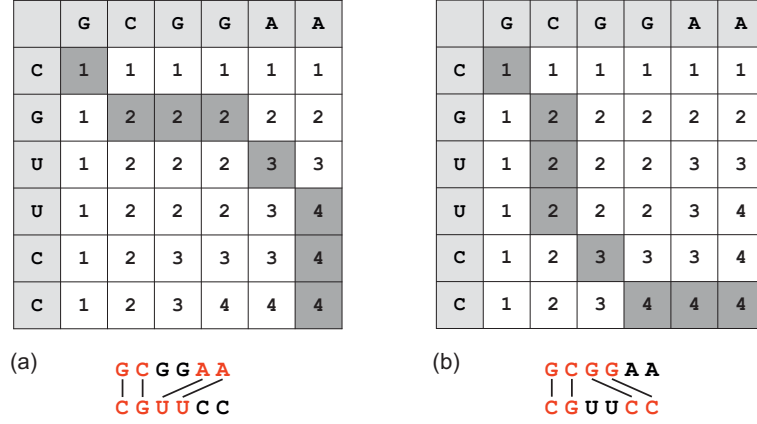


Figure 6: (Color online) Structure recovery algorithm for linear chains.

- For $F_{p-1,q-1}$ consider the corresponding matrix Q (Eq.(23)), chose the maximal element, Q_{\max} of this matrix and compare it with the value $F = F_0 - (f_{p+1,m}^{(1)} + f_{q+1,n}^{(2)} + \tilde{\eta}_{p,q})$; $F_0 = F_{m,n}$ (on the next step we use F instead of F_0). Now,
 - If $Q_{\max} = F$, we look for the next pair (s, r) of linked nucleotides and proceed analogously;
 - If $Q_{\max} < F$, then (according to Eq.(22)) there are on any more pairs of linked nucleotides in the considered branching structure.
- Knowing pairs of linked nucleotides, for example, (p, q) and (s, r) , we reconstruct the structure of the loops between the paired nucleotides by the corresponding statistical weights $f_{p,s}^{(1)}$ and $f_{q,r}^{(2)}$.

The sequence of operations for the structure recovery of RNA-like chains is schematically depicted in the figure 7.

Below we demonstrate on simple example how this algorithm works. Take two sequences $S1$ and $S2$ as shown in Fig.8. The corresponding incidence matrices η' (for intra-matching $S1-S1$), η'' (for intra-matching $S2-S2$), and η (for inter-matching $S1-S2$) are shown in Fig.8 (a), (b) and (c) correspondingly.

The matrices of effective statistical weights $f^{(1)}$ and $f^{(2)}$ of first and second sequences, as well as the ground-state free energy matrix F , are shown in the Fig.9 (a), (b) and (c). The elements $f_{m+1,j}$ and $f_{n+1,j}$, which formally present in the computations, are set to zero: $f_{m+1,j} = f_{n+1,j} = 0$ for all j .

By comparing Fig.9a,b with Fig.9c we see that since $f_{1,7}^{(1)} = f_{1,7}^{(2)} = 2$ and $F_0 = F_{7,7} = 6$, we have $F_{7,7} > f_{1,7}^{(1)} + f_{1,7}^{(2)}$. According to the algorithm described, write the matrix Q corresponding to the element $F_{7,7}$. This matrix Q is depicted in Fig.10a. (Recall that each

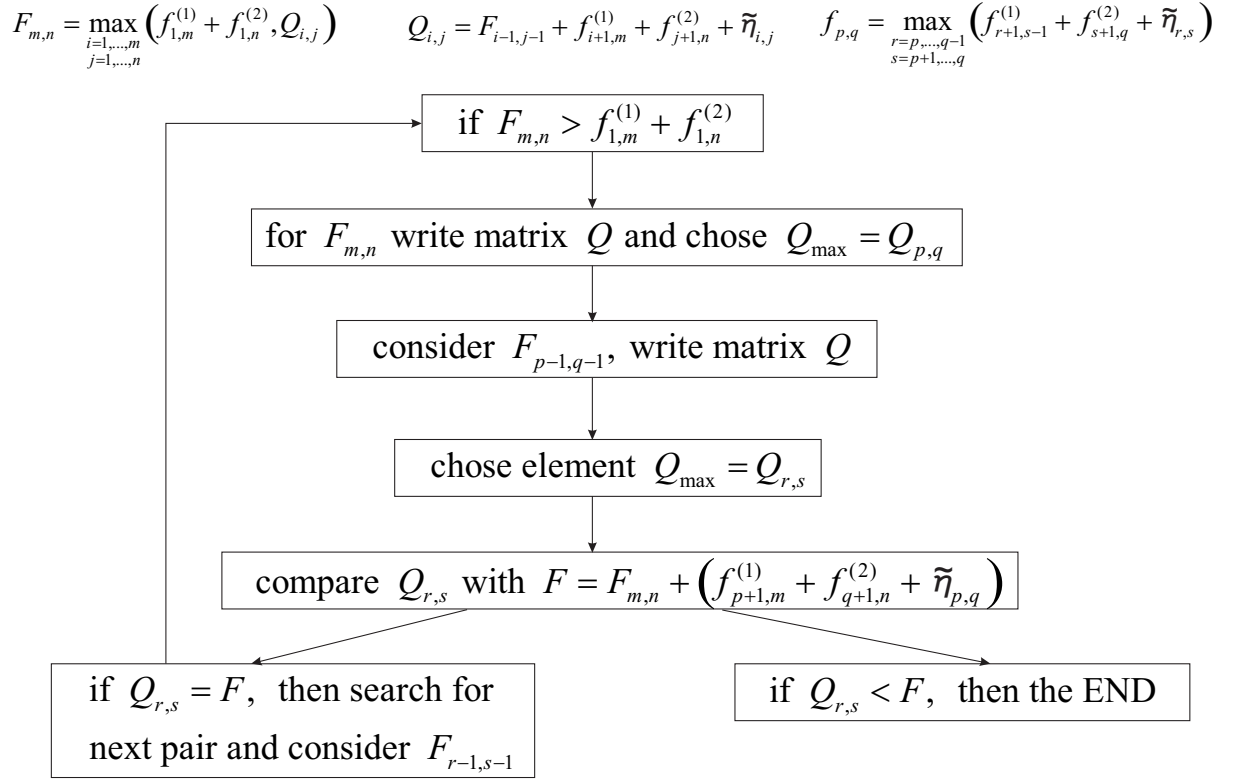


Figure 7: Structure recovery algorithm for RNA molecules.

element $F_{i,j}$ has its own matrix Q of size $i \times j$). We show only those matrices Q which are used for the structure recovery.

The maximal element of the matrix Q depicted in Fig.10a is $Q_{7,7} = 6$, meaning that the 7th nucleotide of S1 interacts with the 7th nucleotide of S2.

To find the next pair of interacting monomers, consider the matrix Q corresponding to the element $F_{6,6}$. This matrix Q is depicted in Fig.10b. It has two maximal elements: $Q_{3,6} = Q_{3,1} = 5$. Thus one has degeneration for the structure under consideration. According to our algorithm, the choice of $Q_{3,1}$ means the interaction of the 3rd nucleotide of the 1st sequence with the 1st nucleotide of the 2nd sequence. At this stage the recovery process is completed. For the choice $Q_{3,6}$ we compute $F^{(1)} = F_0 - (f_{8,7}^{(1)} + f_{8,7}^{(2)} + \tilde{\eta}_{7,7})$. Since $f_{8,7}^{(1)} = f_{8,7}^{(2)} = 0$ and $\tilde{\eta}_{7,7} = 1$, we see that $Q_{3,6} = F^{(1)}$. This means that the 3rd monomer of S1 and the 6th monomer of S2 constitute the next interacting pair. Now we consider $F_{2,5}$. The corresponding matrix Q is shown in the Fig.10c. We see that $Q_{\max} = Q_{2,4} = 2$; $f_{4,6}^{(1)} = 1$; $f_{7,6}^{(2)} = 0$; $\tilde{\eta}_{3,6} = 1$. Since, as before, $F^{(2)} = F^{(1)} - (f_{4,6}^{(1)} + f_{7,6}^{(2)} + \tilde{\eta}_{3,6})$, we see that $Q_{2,4} < F^{(2)}$. Thus, the 2nd and 4th nucleotides do not interact and in the structure there are no more interacting nucleotides. The loop structures can be reconstructed by corresponding statistical weights – see Fig.11.

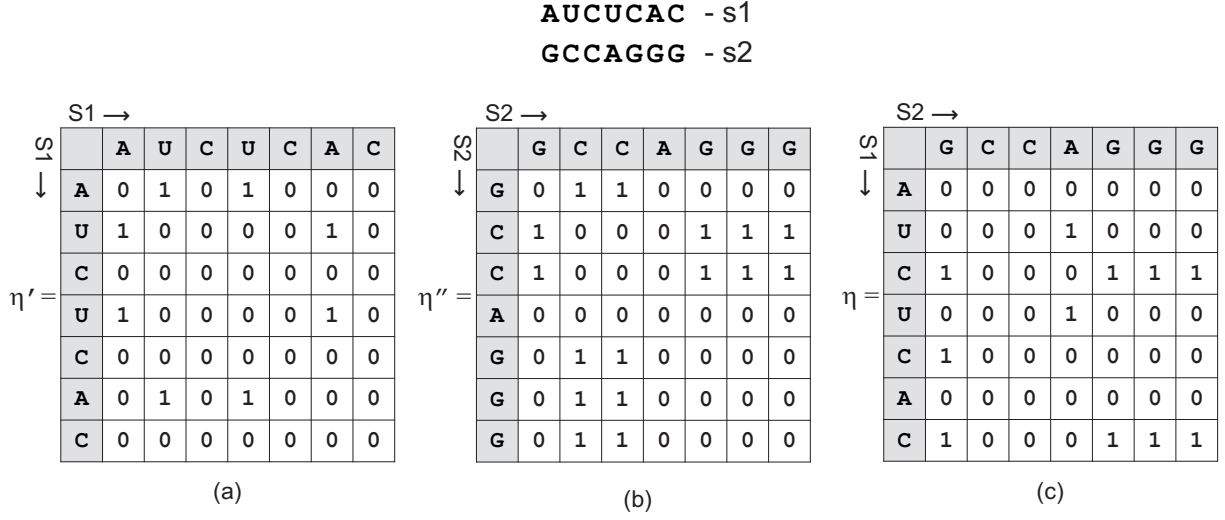


Figure 8: Incidence matrices for pairs of chains with possible clover-leaf structures inside each sequence: (a) intra-matching S1-S1; (b) intra-matching S2-S2; (c) inter-matching S1-S2.

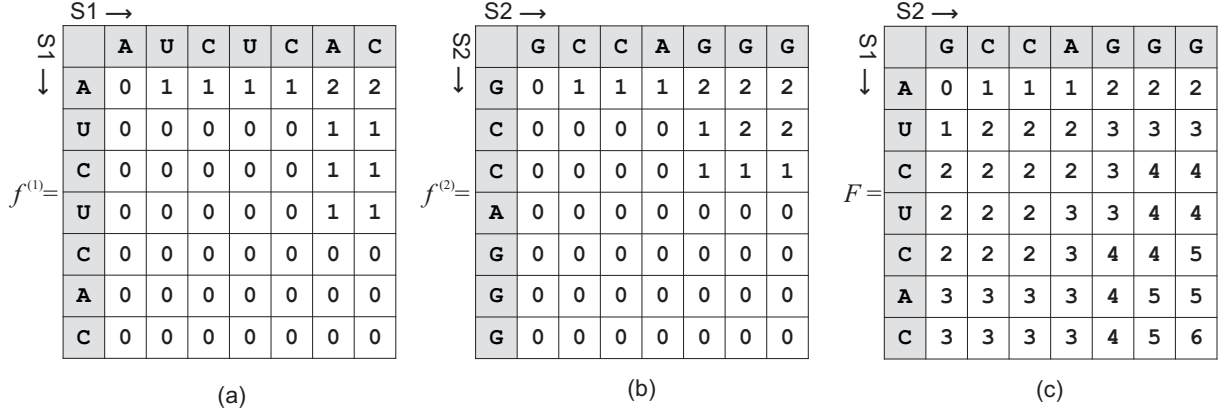


Figure 9: (Color online) Algorithm description: Energies corresponding to incidence matrices in Fig.8: Statistical weights of the 1st (a) and 2nd (b) sequences; (c) Ground-state free energy matrix.

The proposed algorithm is applied to the longer trial sequences shown in Fig.4. Namely, we have performed the structure recovery for three different cases: for linear chains (a) (for them we use the algorithm described in the part 1), for cactus-like chains (b) and for cactus-like chains with the restriction on the size of the minimal loop length (c) (there are no loops less than 4 nucleotides). These structures are depicted in the figures Fig.12a,b and c correspondingly.

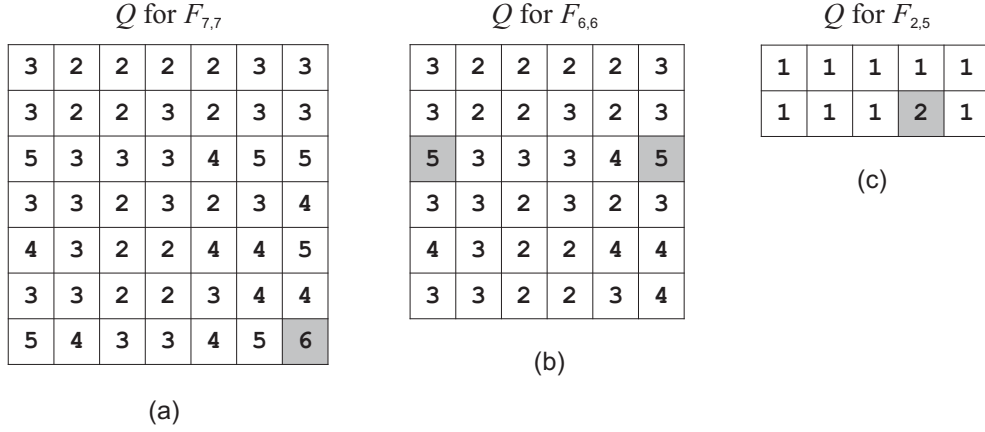


Figure 10: (Color online) Algorithm description: Matrices Q corresponding to: a) $F_{7,7}$; b) $F_{6,6}$; c) $F_{2,5}$.



Figure 11: Algorithm description: Structures recovered from the pair of short sequences shown in Fig.8.

V. CONCLUSION

In this paper we have developed and implemented a new statistical algorithm for quantitative determination of the binding free energy of two heteropolymer sequences under the supposition that each sequence can form a hierarchical cactus-like secondary structure, typical for RNA molecules. For the sequences of lengths m and n the search algorithm is completed in time $\sim O(m^2 \times n^2)$.

We have offered in Section III a constructive way to build a “cost function” characterizing the matching of two *noncoding RNAs* with arbitrary primary sequences. Since base-pairing of two ncRNAs or between ncRNA and DNA plays very important biological role, it is worth estimating theoretically the binding free energy of the ncRNA–target RNA complex by knowing the primary sequences of chains under consideration. Note, that this problem differs from the complete alignment of two RNA sequences: in ncRNA case we align only the sequences of nucleotides which constitute pairs between two RNAs, while the secondary structure of each RNA comes into play only by the combinatorial factors affecting the entropic contribution of chains to the total cost function.

The proposed algorithm is based on two facts: i) the standard alignment problem can be reformulated as a zero-temperature limit of more general statistical problem of binding of two associating heteropolymer chains; ii) the last problem can be straightforwardly generalized onto the sequences with hierarchical cactus-like structures (i.e. of RNA-type). Taking

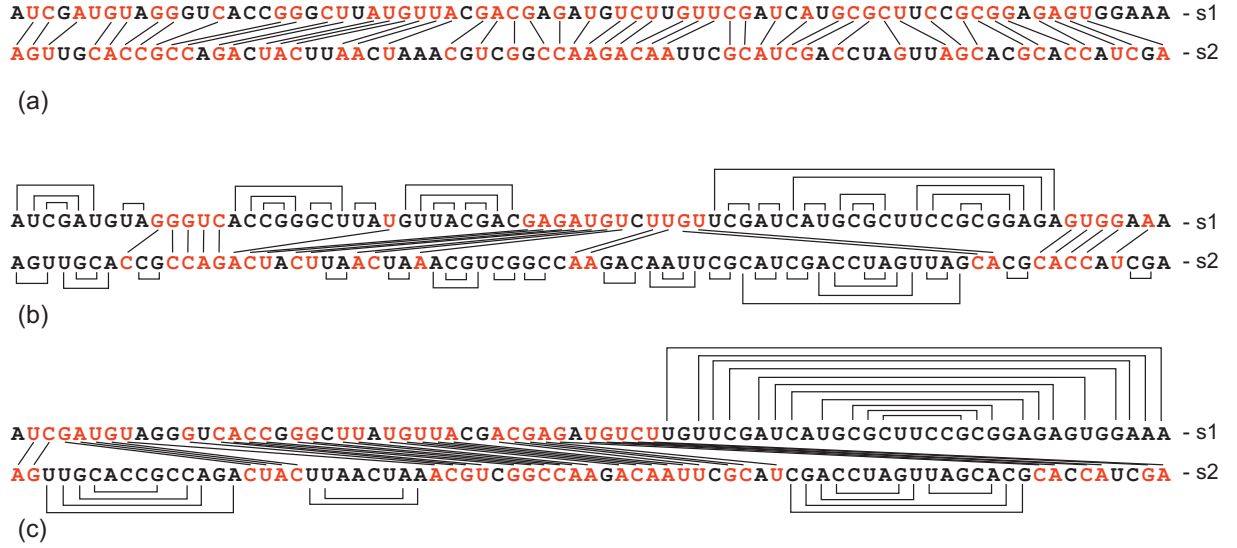


Figure 12: (Color online) Structures recovered from the pair of sequences shown in Fig.4: (a) linear structure; (b) branching structure; (c) branching structure with the restriction on the size of the minimal loop (there are no loops less than 4 nucleotides).

zero-temperature limit at the very end we arrive at the desired ground state free energy with account for entropy of side cactus-like loops.

In this paper we have also demonstrated in detail (see Section IV) how our algorithm enables to solve the *structure recovery* problem, which is in some sense, "inverse" with respect to finding the best matching of two ncRNAs. In particular, we can predict in zero-temperature limit the cactus-like (i.e. the secondary) structure of each ncRNA by knowing only their primary sequences.

In addition we have performed the statistical analysis of a pair of linear and RNA-type random sequences. To avoid the congestion of the paper by the details of computations we have presented these results in Appendix B.

Acknowledgments

We are very grateful to A.A. Mironov for opening for us the world of ncRNAs and to V.A. Avetisov for numerous encouraging discussions concerning the biophysical and statistical aspects of the problem. This work has been partially by the grant ERARSysBio+ #66; M.V. Tamm and O.V. Valba acknowledge the warm hospitality of LPTMS where this work has been completed.

Appendix A: Temperature dependence of the free energy

Analyzing the temperature dependence of the free energy, $F(T)$ for linear and cactus-like chains and have found some significant differences. The figure 13 demonstrates the $F(T)$ –dependencies of the trial sequences shown in Fig.4 under the condition that they form linear or hierarchical cactus-like structures (with loop factor for ideal chains, $\alpha = 3/2$, and with the restrictions on the minimal length, ℓ , of the loop).

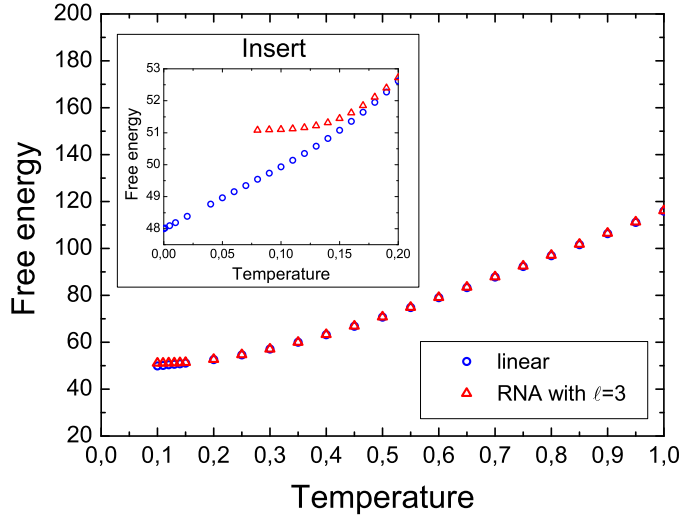


Figure 13: Temperature dependence of the free energy of random trial sequence for linear structures (O) and RNA-like structures (\triangle – with $\ell > 3$).

At high temperatures the $F(T)$ –dependencies for linear and cactus-like (with the minimal loop’s length $\ell = 3$) structures are almost identical. This signals that the creation of any loop of length $\ell > 3$ becomes entropically unfavorable.

At sufficiently low T the $F(T)$ –dependencies for linear and cactus like structures (with $\ell > 3$) deviate from each other. This deviation has rather transparent physical explanation. Represent the free energy $F(T)$ at $T \rightarrow 0$ in the following generic form

$$F = F_0 + T \ln W + T e^{-u/T} \quad (\text{A1})$$

where F_0 is the ground state energy and W is the number of states with the same energy (degeneration). According to (A1) the slope of the curve $F(T)$ at $T \rightarrow 0$ determines the degree of the degeneration. Decrease of the slope for cactus-like structures indicates that the creation of hierarchical “cactuses” and account for entropy of loops removes the degeneration.

Appendix B: Statistical analysis of a pair of random sequences

We have analyzed the basic statistical properties of a pair of random sequences. For simplicity, we considered the chains of the same length n . It has been shown in [29] that for *linear sequences* the ground state free energy in the so-called "Bernoulli matching approximation" has the following behavior at $n \gg 1$:

$$\begin{aligned}\langle F \rangle &\approx \frac{2}{1 + \sqrt{c}} n + f(c) \langle \chi \rangle n^{1/3} \\ \sigma &\equiv \sqrt{\langle F^2 \rangle - \langle F \rangle^2} \approx \sqrt{\langle \chi^2 \rangle - \langle \chi \rangle^2} f(c) n^{1/3}\end{aligned}\tag{B1}$$

where $f(c) = \frac{c^{1/6}(\sqrt{c}-1)^{1/3}}{\sqrt{c}+1}$ (see [29] for details), c is the number of different nucleotides (in our case $c = 4$) and χ is some random variable with known n -independent distribution ($\langle \chi \rangle = -1.7711\dots$ and $\langle \chi^2 \rangle - \langle \chi \rangle^2 = 0.8132\dots$).

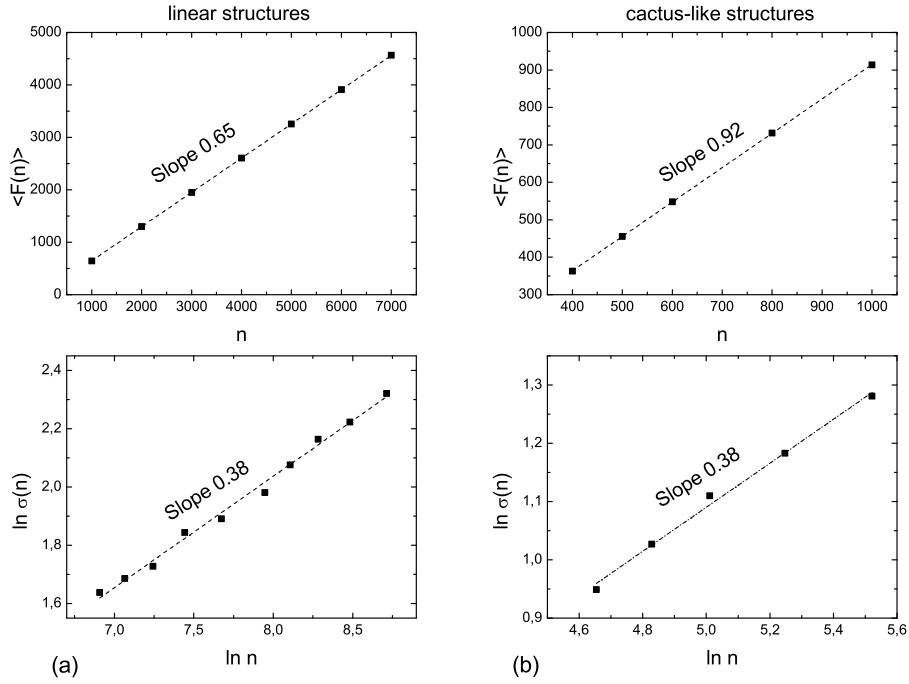


Figure 14: Plots of the average free energy, $\langle F(n) \rangle$ (linear scale), and its fluctuations, $\sigma(n)$ (double logarithmic scale) in zero-temperature limit for: a) linear chains and b) cactus-like chains.

In the Fig.14a we have plotted $\langle F(n) \rangle$ (in the linear scale) and $\sigma(n)$ (in the double logarithmic scale). One sees that the slope $k_1 \approx 0.65$ in Fig.14a is in very good agreement with the value $k_1 = \lim_{n \rightarrow \infty} \frac{\langle F \rangle}{n} \rightarrow \frac{2}{3}$ computed from the 1st of equations (B1), while the slope 0.38 in the Fig.14b is close to the exponent $\frac{1}{3}$ in the 2nd line of (B1). The averaging has been performed over 200 different randomly chosen structures with uniform distribution of $c = 4$ nucleotides.

The similar analysis have been performed for sequences with the possibility of cactus-like structure formation. The plots of $\langle F(n) \rangle$ and $\sigma(n)$ are shown in the figures 14c (in linear scale and in double logarithmic scale correspondingly). One sees that again, as for linear sequences, $\langle F(n) \rangle = k_c n$ for large n , but the coefficient $k_c \approx 0.92$ is larger than k_l what signals the large number of pairs in the ground state, leading to the loop creation. The slope in the Fig.14d allows one to conclude that the loop creation does not affect the universality class of the fluctuations and it remains the same as for linear sequences.

-
- [1] V. Ambros, *Cell* **107**, 862 (2001)
 - [2] G. Storz, *Science* **296**, 1260 (2002)
 - [3] P. Navarro, S. Pichard, C. Ciaudo, P. Avner, C. Rougeulle, *Genes & Development* **19** 1474 (2005)
 - [4] S.R. Eddy, *Cell* **109**, 137 (2002)
 - [5] V. Pande, A. Grosberg, T. Tanaka, *Rev. Mod. Phys.* **72**, 259 (2000)
 - [6] E. Rivas, S.R. Eddy, *J. Mol. Biol.* **285**, 205 (1999)
 - [7] S.B. Needleman and C.D. Wunsch, *J. Mol. Biol.* **48**, 443 (1970)
 - [8] T.F. Smith and M.S. Waterman, *J. Mol. Biol.* **147**, 195 (1981); *Adv. Appl. math.* **2**, 482 (1981)
 - [9] M.S. Waterman, L. Gordon, and R. Arratia, *Proc. Natl. Acad. Sci. USA*, **84**, 1239 (1987)
 - [10] S.F. Altschul et. al., *J. Mol. Biol.* **215**, 403 (1990)
 - [11] D. Sankoff and J. Kruskal, *Time Warps, String Edits, and Macromolecules: The theory and practice of sequence comparison* (Addison Wesley, Reading, Massachusetts, 1983)
 - [12] A. Apostolico and C. Guerra, *Algoritmica*, **2**, 315 (1987)
 - [13] R. Wagner and M. Fisher, *J. Assoc. Comput. Mach.* **21**, 168 (1974)
 - [14] D. Gusfield, *Algorithms on Strings, Trees, and Sequences* (Cambridge University Press, Cambridge, 1997)
 - [15] J. Boutet de Monvel, *European Phys. J. B* **7**, 293 (1999); *Phys. Rev. E* **62**, 204 (2000)
 - [16] V. Chvátal and D. Sankoff, *J. Appl. Probab.* **12**, 306 (1975)
 - [17] J. Deken, *Discrete Math.* **26**, 17 (1979)
 - [18] J.M. Steele, *SIAM J. Appl. Math.* **42**, 731 (1982)
 - [19] V. Dancik and M. Paterson, in *STACS94, Lecture Notes in Computer Science*, **775**, 306 (Springer: New York, 1994)
 - [20] K.S. Alexander, *Ann. Appl. Probab.* **4**, 1074 (1994)
 - [21] M. Kiwi, M. Loeb, and J. Matousek, in *Lecture Notes in Computer Science*, **2976** 302 (Springer: Berlin, 2004)
 - [22] M. Zhang and T. Marr, *J. Theor. Biol.* **174**, 119 (1995)
 - [23] T. Hwa and M. Lassig, *Phys. Rev. Lett.* **76**, 2591 (1996)

- [24] R. Bundschuh, Eur. Phys. J. B **22**, 533 (2001)
- [25] M.S. Waterman, Bull. Math. Biol. **46**, 473 (1984)
- [26] M.S. Waterman and M. Vingron, Statistical Science, **9**, 387 (1994)
- [27] R. Bundschuh, T. Hwa, Discrete Appl. Math. **104**, 113 (2000).
- [28] D. Drasdo, T. Hwa, M. Lassig, J. Comp. Biol. **7**, 115 (2000)
- [29] S.N. Majumdar and S. Nechaev, Phys. Rev. E **69**, 011103 (2004).
- [30] C.A. Tracy and H. Widom, Comm. Math. Phys. **159**, 151 (1994); see also Proc. of ICM, Beijing, Vol. I, 587 (2002).
- [31] For a recent review of the appearance of Tracy–Widom distribution in several physics problems, see S.N. Majumdar (Les Houches lecture notes on ‘Complex Systems’, 2007), arXiv: cond-mat/0701193.
- [32] M.V. Tamm, S.K. Nechaev, Phys. Rev. E **78**, 011903 (2008)
- [33] L. Comtet, *Advanced Combinatorics: The Art of Finite and Infinite Expansions*, (Dordrecht: Reidel, 1974)
- [34] M. Mueller, Phys. Rev. E, **67**, 021914 (2003)
- [35] A.M. Gutin, A.Yu. Grosberg, E.I. Shakhnovich, J.Phys. A: Math. Gen. **26**, 1037 (1993)
- [36] P. de Gennes, Biopolymers, **6**, 715 (1968)
- [37] I.Ya. Erukhimovich, Vysokomolek. Soed., **20B**, 10 (1978) (*in Russian*)
- [38] This general description covers both cases (DNA and RNA) by a straightforward redefinition of letters.
- [39] The cooperativity means that if two links are connected with each other, then the two adjacent links have larger probability to be also connected.